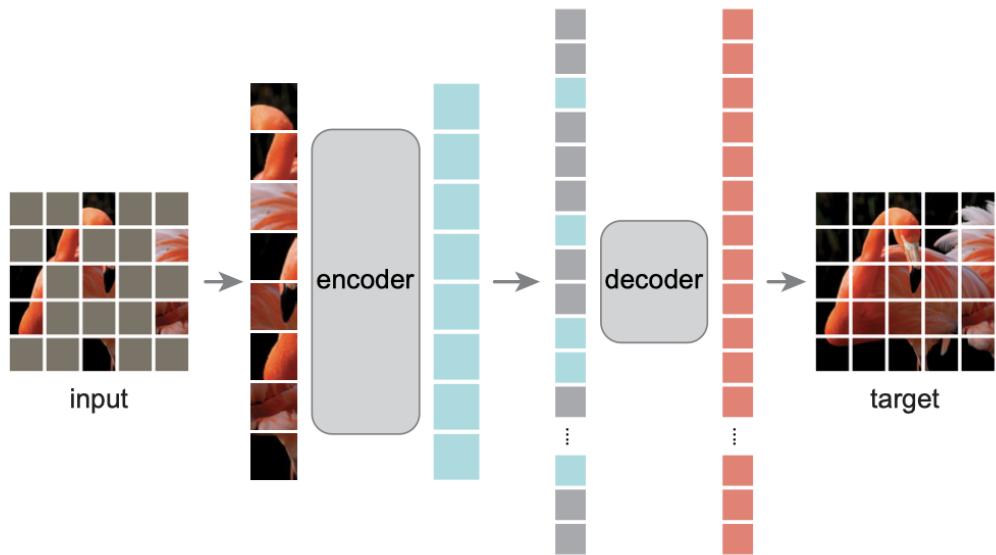


# Generative Graph Self-Supervised Learning

Mar. 8, 2023  
Liang Wang

# Graph MAE Family



## •MAE vs Bert

- **MAE is an AutoEncoder architecture, while Bert is not:** MAE involves pixel reconstruction and requires a decoder, whereas Bert predicts labels and needs a predictor.
- **The mask ratio in MAE is higher** because there is more redundancy in image information.
- **The mask token in MAE is not input into the encoder:** the encoder focuses solely on feature extraction and not reconstruction. In MAE, the encoder acts as a feature extractor, and during the fine-tuning phase, only the encoder is retained. To maintain consistency with the inputs during the fine-tuning phase, the mask token is not input into the encoder).

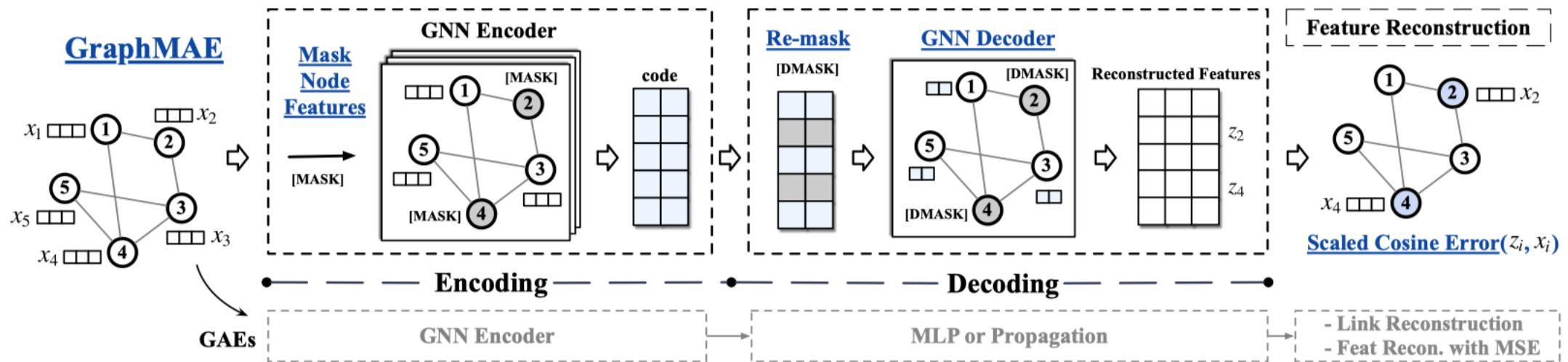
## •MAE vs DAE

- Essentially, it belongs to a type of Denoising AutoEncoder (DAE), but the difference is: **MAE is asymmetric**, where the encoder only inputs the unmasked parts, and the decoder inputs both parts; **whereas DAE is symmetric**.

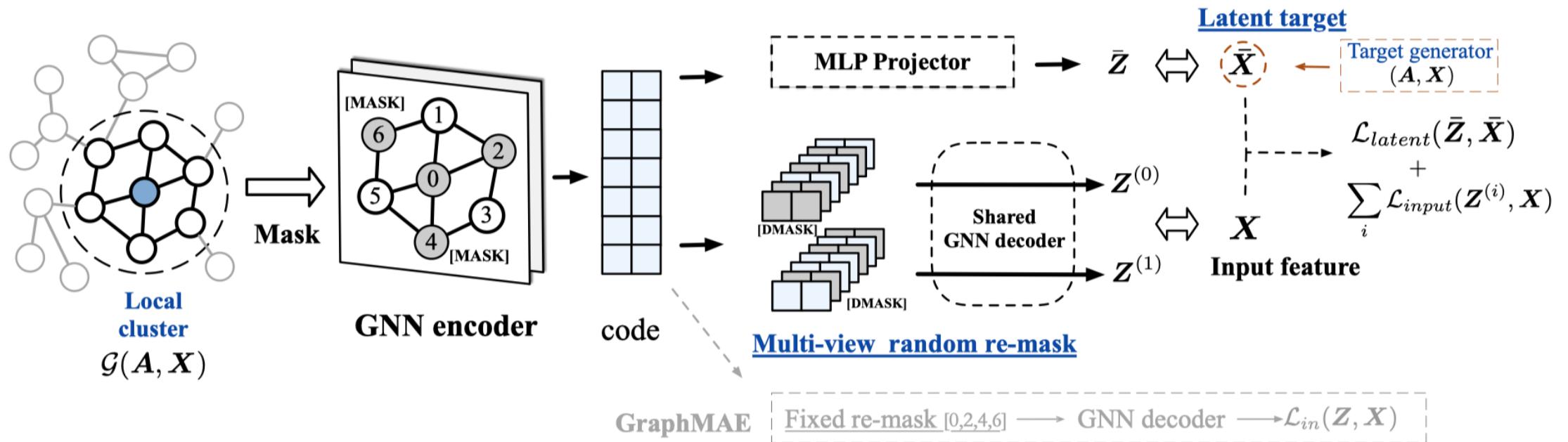
## •Optimization Objective

- MSE on normalized patches is better than MSE**

# GraphMAE, KDD, 2022



# GraphMAE2, WWW, 2023



# MaskGAE, KDD, 2023

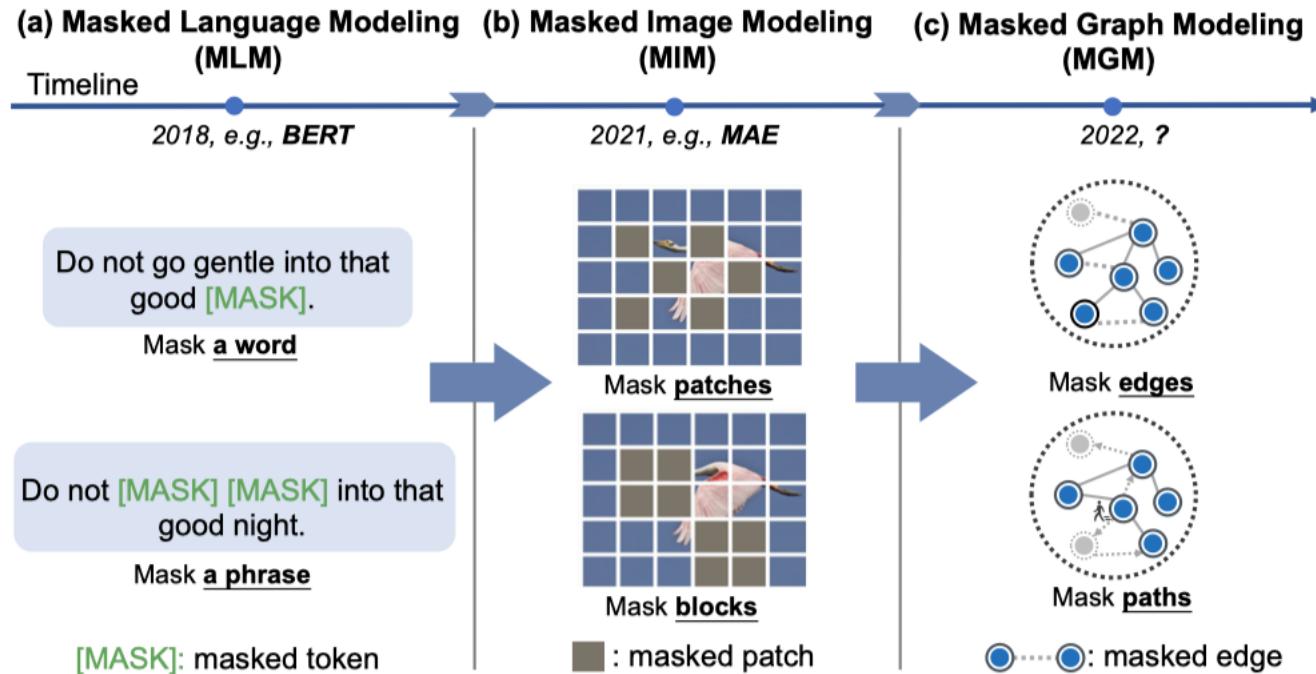


Figure 1: **From left to right:** illustrative examples of (a) masked language modeling (MLM), (b) masked image modeling (MIM), and (c) masked graph modeling (MGM) paradigms with different masking strategies. Similar to MLM and MIM, the goal of MGM is to learn representations by predicting randomly masked edges based on visible structure.

# MaskGAE, arXiv, 2022

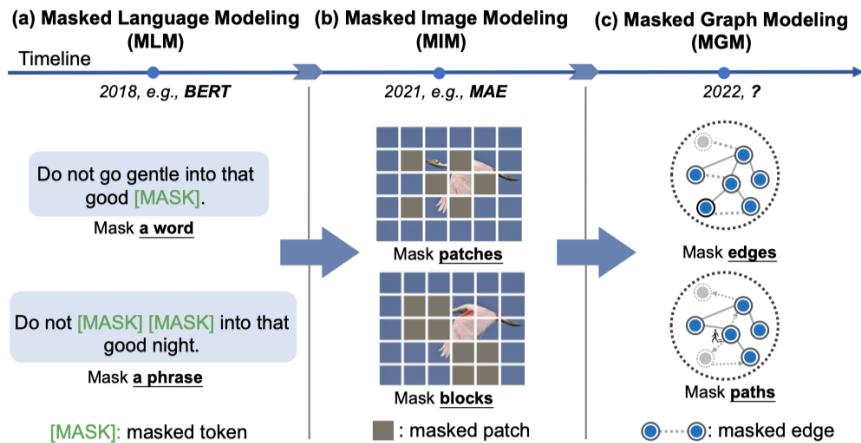


Figure 1: **From left to right:** illustrative examples of (a) masked language modeling (MLM), (b) masked image modeling (MIM), and (c) masked graph modeling (MGM) paradigms with different masking strategies. Similar to MLM and MIM, the goal of MGM is to learn representations by predicting randomly masked edges based on visible structure.

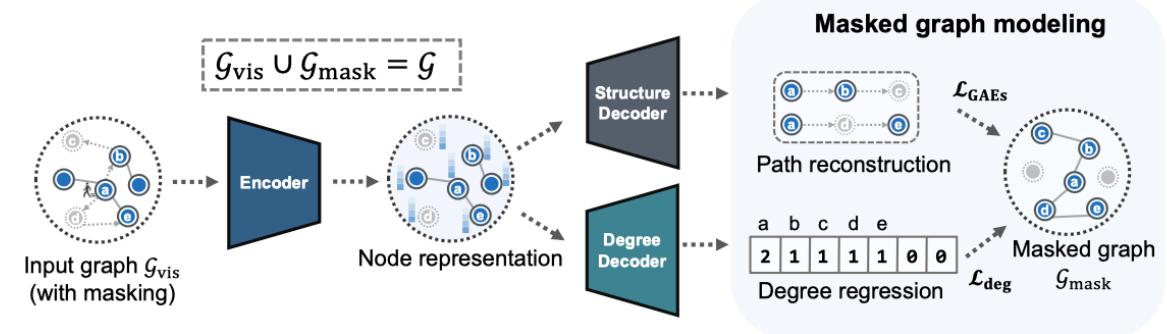
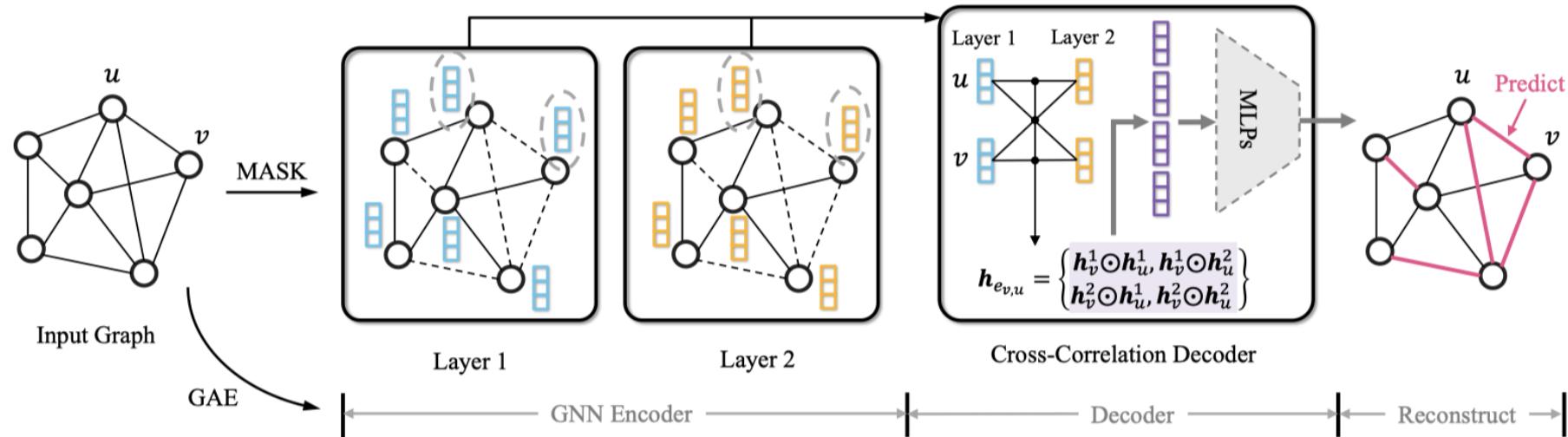


Figure 2: Overview of MaskGAE framework, performing masked graph modeling with an asymmetric encoder-decoder design. During the self-supervised learning phase, an input graph  $\mathcal{G}_{vis}$  is provided, but with some of paths (a set of adjacent edges) are masked. The goal of the model is then to learn to predict the masked edges in the graph  $\mathcal{G}_{mask}$  and the degree of the associated nodes, respectively.

# S2GAE, WSDM, 2023



**Figure 1: The proposed S2GAE architecture. Given a graph, we first apply direction-aware graph masking strategies to disturb it, obtaining perturbed graph and masked edge set. Then, the perturbed graph is fed into GNN encoder to produce hidden representations. Next, a tailored cross-correlation decoder is designed to reconstruct these masked edges by capturing the cross-correlation of their end nodes from multi-granularity representations. Finally, the whole framework is trained end-to-end by maximizing the likelihood of the masked edge set.**

- [1] Qiaoyu Tan, Ninghao Liu, Xiao Huang, Rui Chen, Soo-Hyun Choi, Xia Hu. “MGAE: Masked Autoencoders for Self-Supervised Learning on Graphs.” arXiv. January 7, 2022
- [2] Qiaoyu Tan, Ninghao Liu, Xiao Huang, Soo-Hyun Choi, Li Li, Rui Chen, Xia Hu. “S2GAE: Self-Supervised Graph Autoencoders Are Generalizable Learners with Graph Masking.” In WSDM. 2023

# GMAE, arXiv, 2022

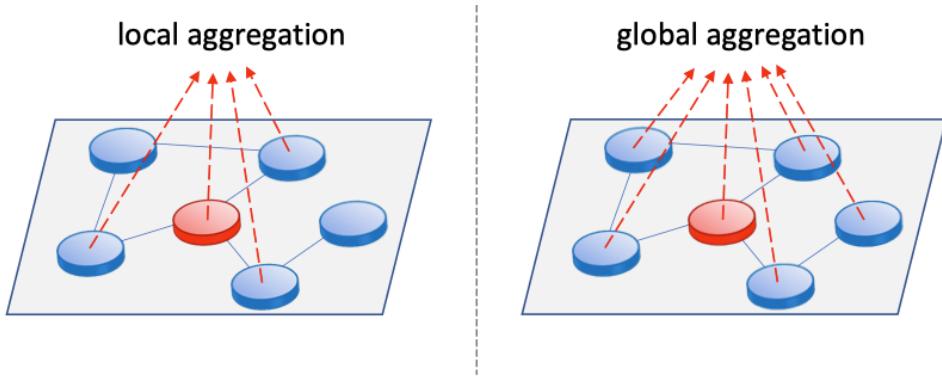


Fig. 1. An illustration of different aggregation mechanisms used by GCNs and transformers. The red node is the target node. The first picture (left) shows the local aggregation used by a GCN layer, where the embedding of the target node is obtained by aggregating its neighbors. The second picture (right) shows the global aggregation used by a transformer layer, where the embedding is obtained by aggregating all the nodes in the graph.

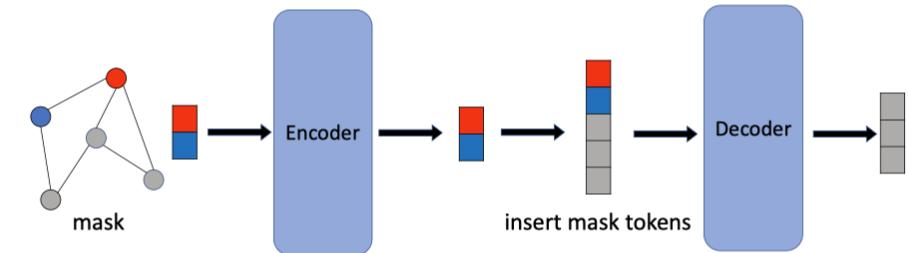


Fig. 2. GMAE framework. We first randomly mask some of the nodes (the grey nodes). The feature matrix and positional embeddings of the non-masked nodes (the red node and the blue node) are fed into an encoder to obtain the representations of the non-masked nodes. Then we insert a shared learnable mask token to represent the masked nodes. The decoder reconstructs the original features of the masked nodes.

# HGMAE, AAAI, 2023

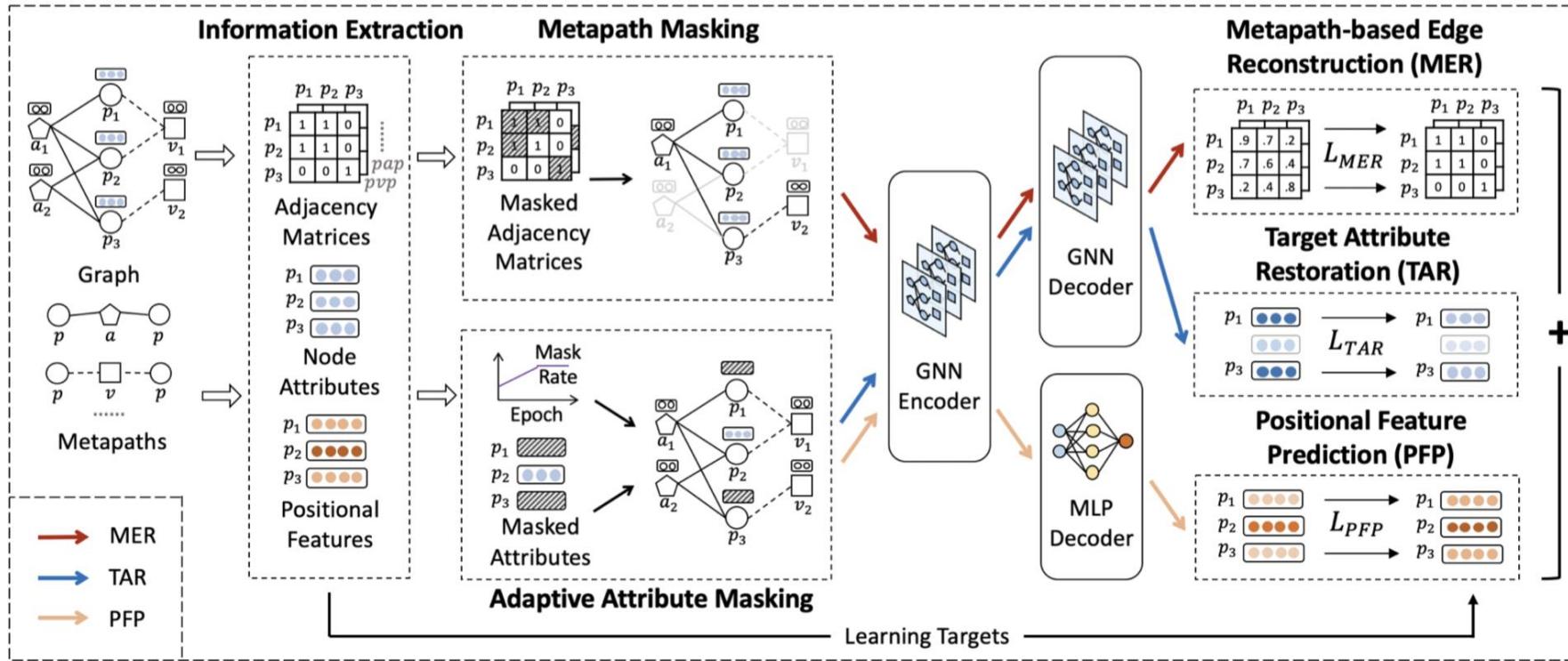


Figure 2: The overall framework of HGMAE: we first extract the metapath-based adjacency matrices, node attributes and positional features from the graph. We then design two masking techniques (i.e., metapath masking and adaptive attribute masking) to mask the inputs. Later, we feed the masked inputs into the encoder and decoders sequentially, and optimize them via three training strategies (i.e., metapath-based edge reconstruction, target attribute restoration, and positional feature prediction). The proposed strategies enable the model to capture comprehensive graph information and address identified challenges.

# BatmanNet, arXiv, 2022

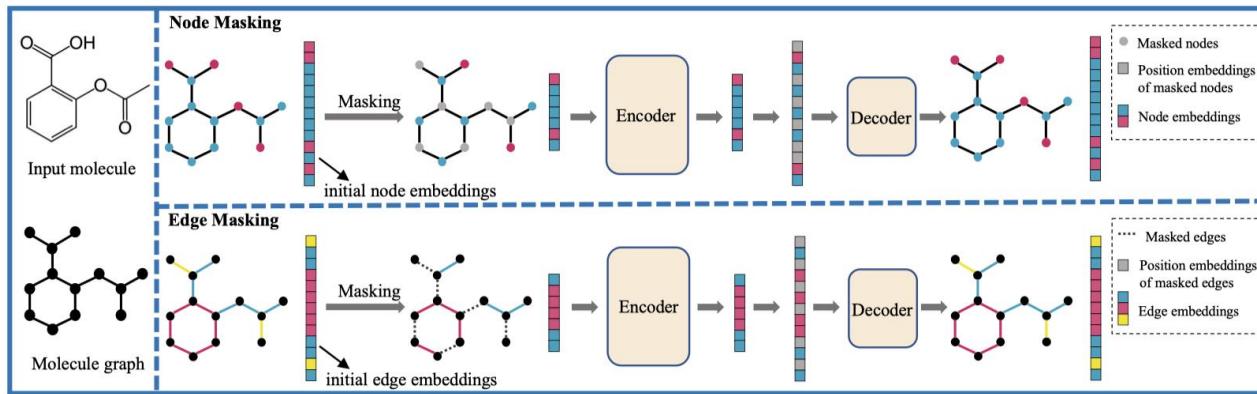
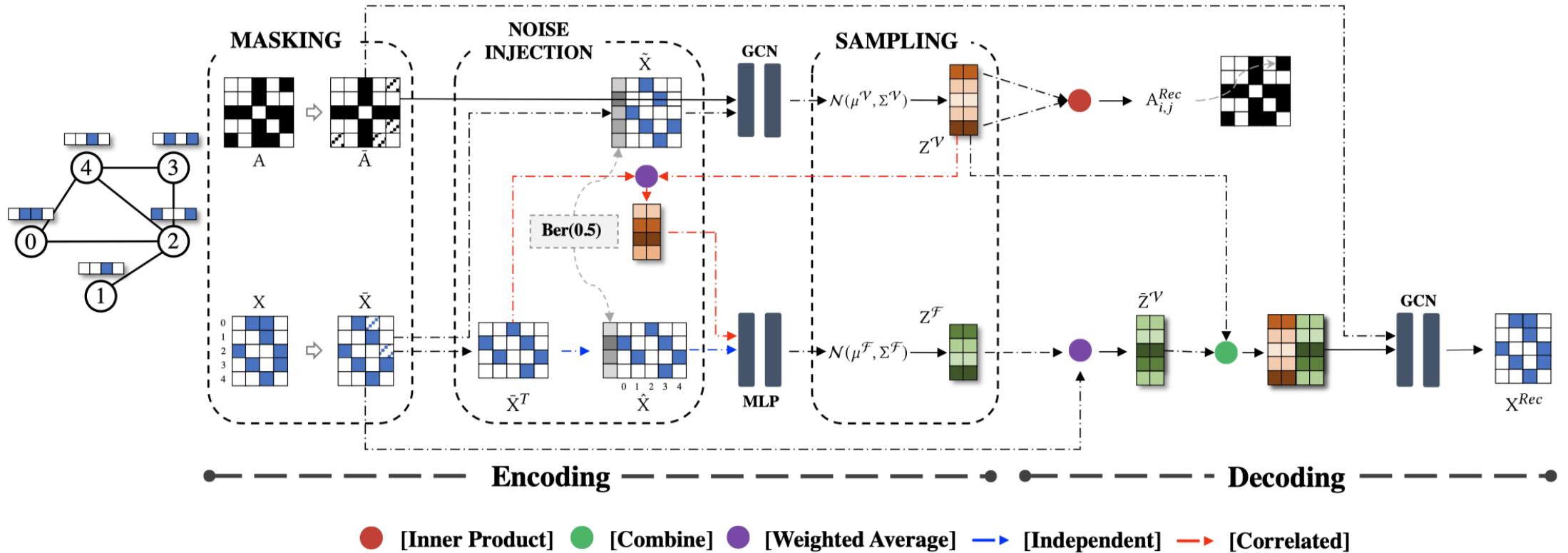


Figure 1: Illustration of the designed self-supervised tasks of BatmanNet. A very high portion of nodes or edges is randomly masked, then the BatmanNet is pre-trained to reconstruct the original molecule from the latent representation and mask tokens.

Dataset	BBBP	SIDER	ClinTox	BACE	Tox21	ToxCast	MUV	HIV
#Molecules	2039	1427	1478	1513	7831	8575	93087	41127
#tasks	1	27	2	1	12	617	17	1
TF_Robust	0.860 <sub>(0.087)</sub>	0.607 <sub>(0.033)</sub>	0.765 <sub>(0.085)</sub>	0.824 <sub>(0.022)</sub>	0.698 <sub>(0.012)</sub>	0.585 <sub>(0.031)</sub>	/	/
GraphConv	0.877 <sub>(0.036)</sub>	0.593 <sub>(0.035)</sub>	0.845 <sub>(0.051)</sub>	0.854 <sub>(0.011)</sub>	0.772 <sub>(0.041)</sub>	0.650 <sub>(0.025)</sub>	/	/
Weave	0.837 <sub>(0.065)</sub>	0.543 <sub>(0.034)</sub>	0.823 <sub>(0.023)</sub>	0.791 <sub>(0.008)</sub>	0.741 <sub>(0.044)</sub>	0.678 <sub>(0.024)</sub>	/	/
SchNet	0.847 <sub>(0.024)</sub>	0.545 <sub>(0.038)</sub>	0.717 <sub>(0.042)</sub>	0.750 <sub>(0.033)</sub>	0.767 <sub>(0.025)</sub>	0.679 <sub>(0.021)</sub>	/	/
MGCN	0.850 <sub>(0.064)</sub>	0.552 <sub>(0.018)</sub>	0.634 <sub>(0.042)</sub>	0.734 <sub>(0.030)</sub>	0.707 <sub>(0.016)</sub>	0.663 <sub>(0.009)</sub>	/	/
N_GRAM	0.912 <sub>(0.013)</sub>	0.632 <sub>(0.005)</sub>	0.855 <sub>(0.037)</sub>	0.876 <sub>(0.035)</sub>	0.769 <sub>(0.027)</sub>	/	/	/
PretrainGNN	0.915 <sub>(0.040)</sub>	0.614 <sub>(0.006)</sub>	0.762 <sub>(0.058)</sub>	0.851 <sub>(0.027)</sub>	0.811 <sub>(0.015)</sub>	0.714 <sub>(0.019)</sub>	<b>0.839</b> <sub>(0.012)</sub>	0.814 <sub>(0.003)</sub>
GraphMAE	0.896 <sub>(0.007)</sub>	0.652 <sub>(0.001)</sub>	0.850 <sub>(0.007)</sub>	0.863 <sub>(0.002)</sub>	0.794 <sub>(0.003)</sub>	0.679 <sub>(0.005)</sub>	0.796 <sub>(0.008)</sub>	0.790 <sub>(0.004)</sub>
GROVER <sub>base</sub>	0.936 <sub>(0.008)</sub>	0.656 <sub>(0.006)</sub>	0.925 <sub>(0.013)</sub>	0.878 <sub>(0.016)</sub>	0.819 <sub>(0.020)</sub>	0.723 <sub>(0.010)</sub>	0.757 <sub>(0.038)</sub>	0.764 <sub>(0.029)</sub>
GROVER <sub>large</sub>	0.940 <sub>(0.019)</sub>	0.658 <sub>(0.023)</sub>	<b>0.944</b> <sub>(0.021)</sub>	0.894 <sub>(0.028)</sub>	0.831 <sub>(0.025)</sub>	<b>0.737</b> <sub>(0.010)</sub>	0.765 <sub>(0.048)</sub>	0.777 <sub>(0.036)</sub>
BatmanNet	<b>0.944</b> <sub>(0.007)</sub>	<b>0.665</b> <sub>(0.004)</sub>	0.889 <sub>(0.012)</sub>	<b>0.911</b> <sub>(0.015)</sub>	<b>0.832</b> <sub>(0.013)</sub>	0.733 <sub>(0.009)</sub>	0.786 <sub>(0.023)</sub>	<b>0.820</b> <sub>(0.053)</sub>

# SeeGera, WWW, 2023



# WGDN, AAAI, 2023

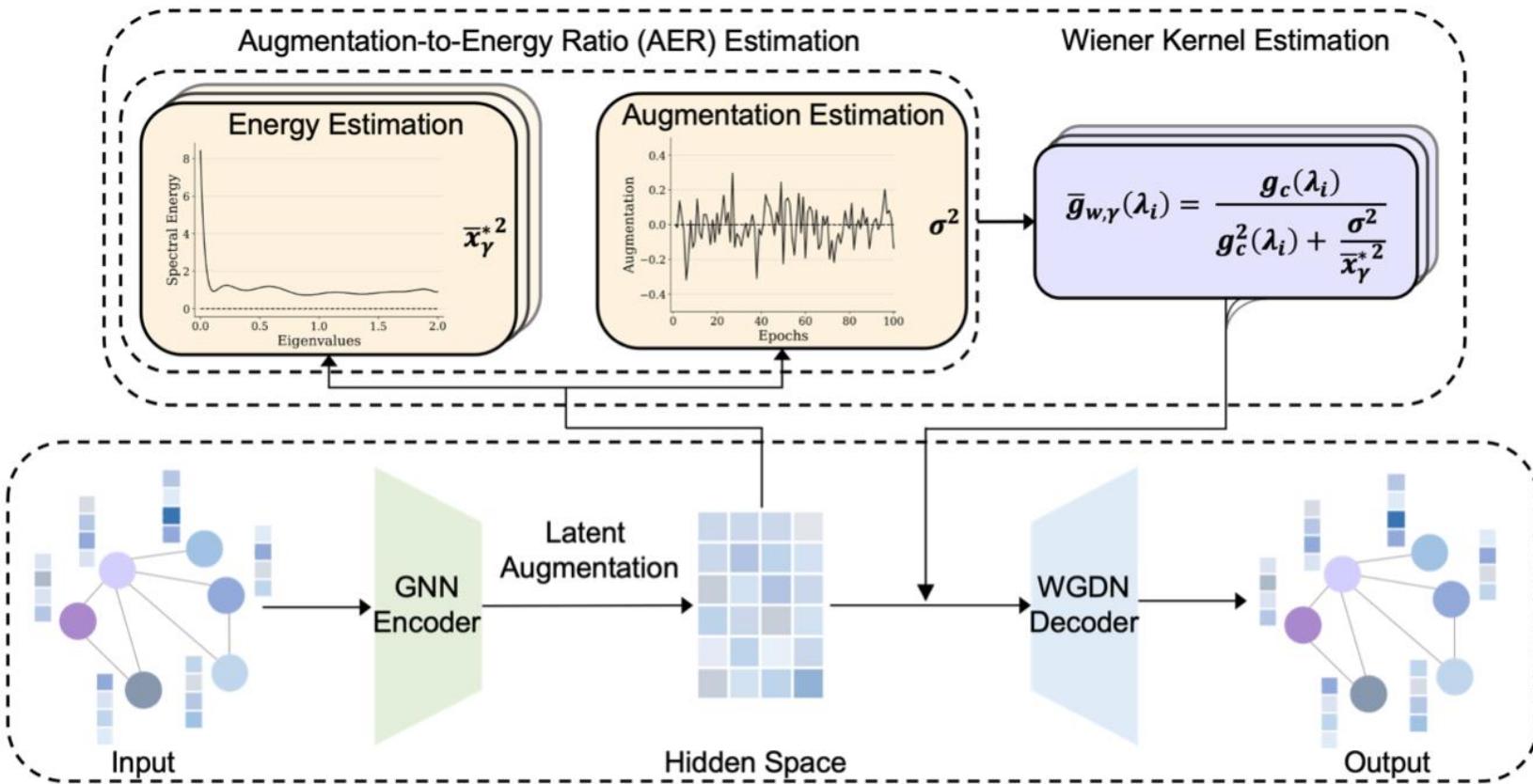


Figure 2: The autoencoder framework of WGDN for graph SSL. Given the augmented latent representations, graph wiener filter is approximated via estimating spectral energy and augmentations adaptively. With such, WGDN permits the stable feature reconstruction from the augmented latent space for representation learning.

# MoCE, ICLR, 2023

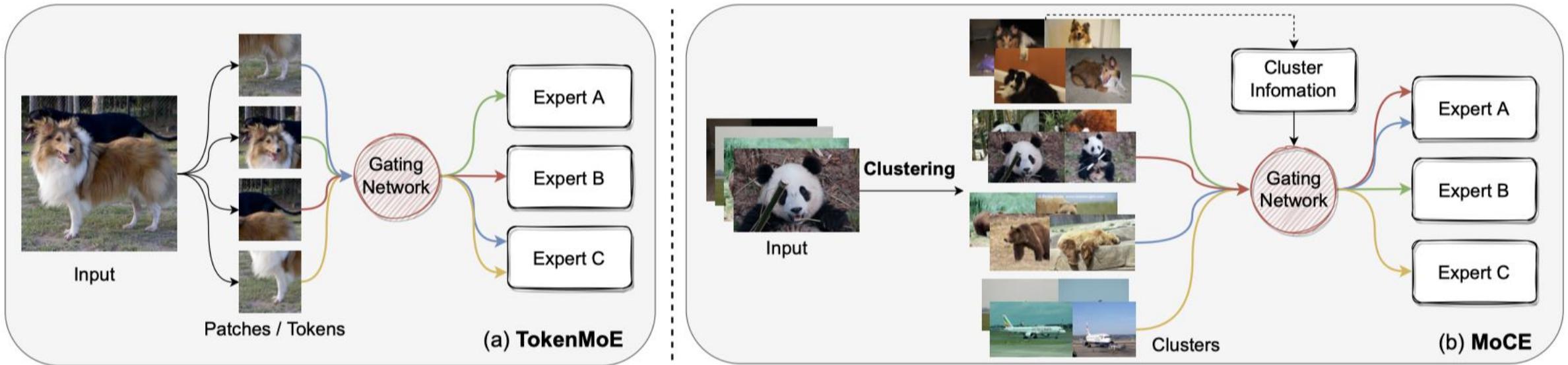
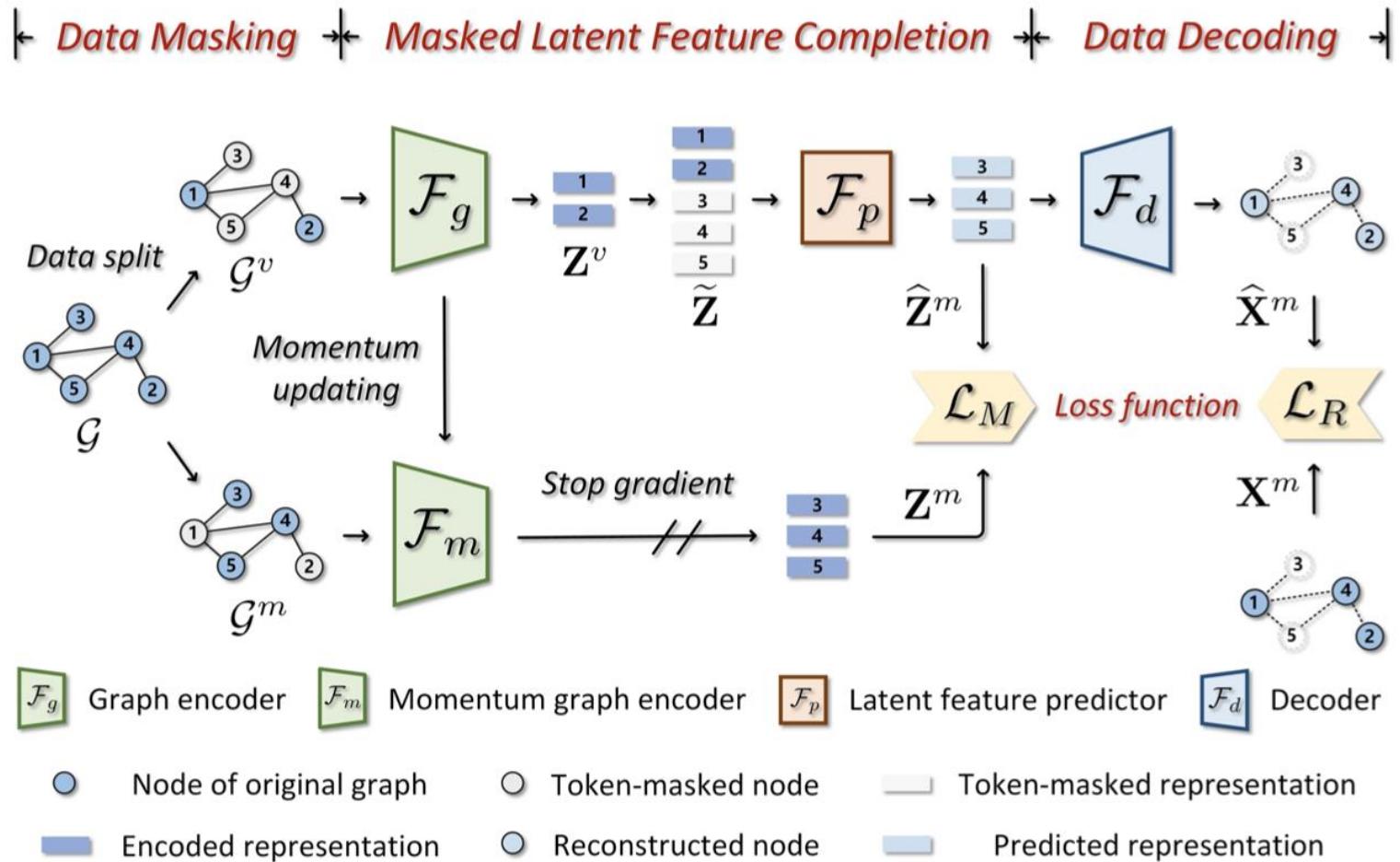


Figure 2: Model design comparison between (a) TokenMoE (Riquelme et al., 2021) and (b) MoCE. Both methods utilize the multi-expert architecture with the main difference about the input of the gating network. MoCE adopts the corresponding cluster embedding of the current token as in Eqn. 4, instead of the token embedding in Eqn. 3.2. Therefore, each expert can be trained by semantically similar images to alleviate the negative transfer phenomenon.

# RARE, arXiv, 2023



# Theoretical Analysis

# Understand MAE, arXiv, 2022

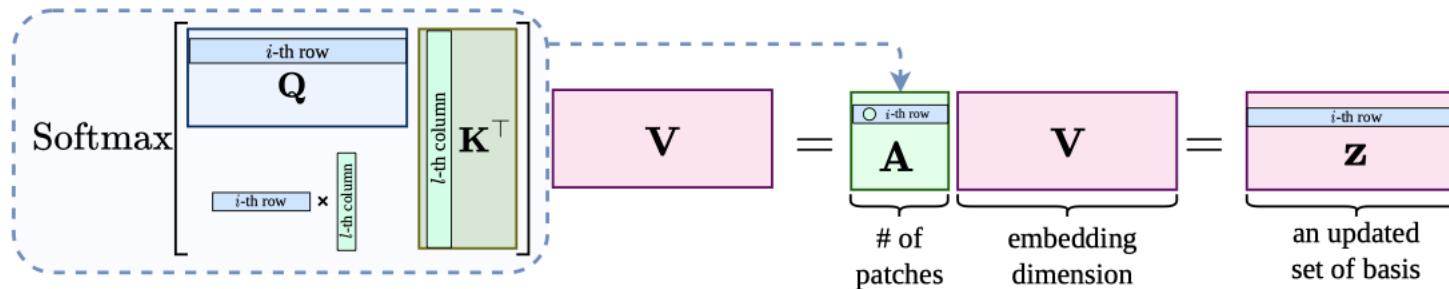


Figure 2: The  $i$ -th patch embedding from the scaled dot-product attention is a convex combination of the attention weights, which encodes the interactions among all patch embeddings. Best viewed in color.

# U-MAE, NeurIPS, 2022

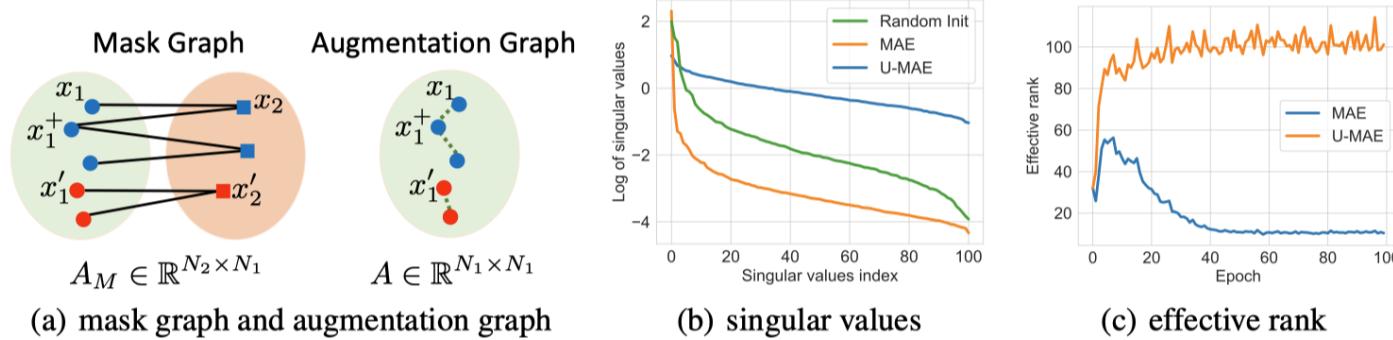


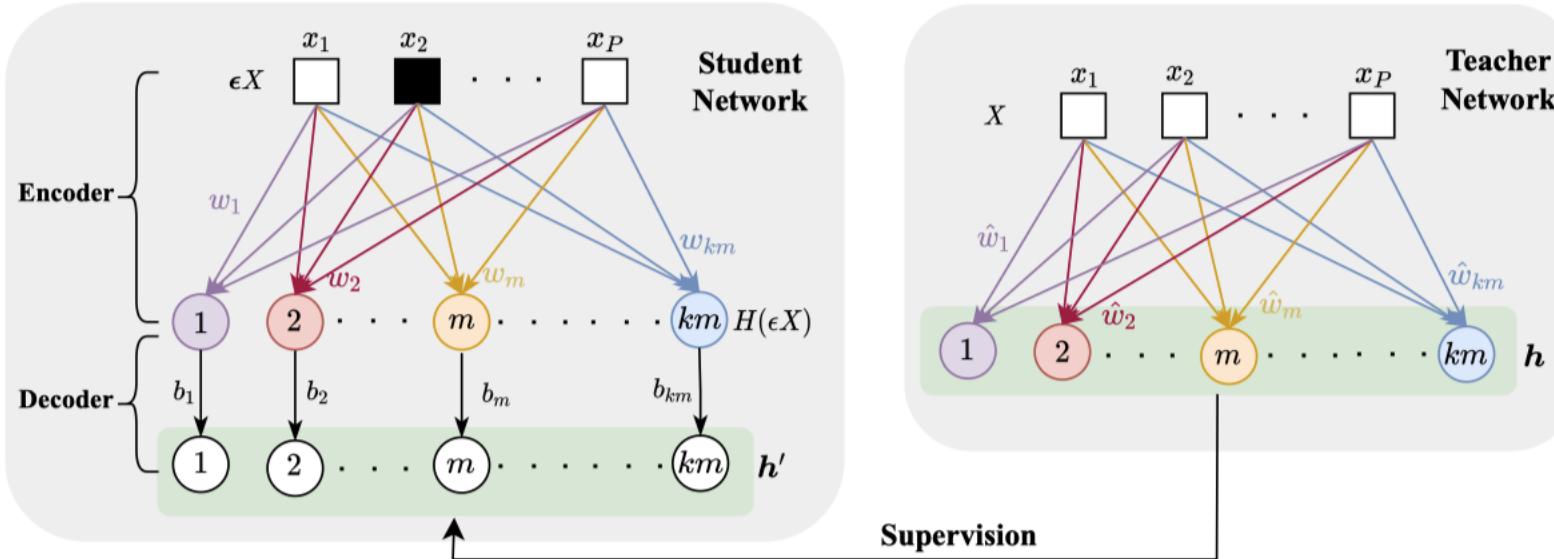
Figure 1: (a) An illustration of the mask graph and the corresponding augmentation graph of MAE. Different colors denote different belonging classes. (b) Comparison of the singular values of learned features with 1) random initialization, 2) MAE loss, 3) our U-MAE loss. (c) The changing process of effective rank [24] of the encoded features trained with different objectives (MAE and U-MAE).

**Theorem 3.4.** *Under Assumption 3.1, MAE's reconstruction loss (Eq. (2)) can be lower bounded by the alignment loss between positive pairs  $(x_1, x_1^+) \sim \mathcal{A}(x_1, x_1^+)$ ,*

$$\mathcal{L}_{\text{MAE}}(h) \geq \frac{1}{2}\mathcal{L}_{\text{align}}(h) - \varepsilon + \text{const} = -\frac{1}{2}\mathbb{E}_{x_1, x_1^+} h(x_1)^\top h(x_1^+) - \varepsilon + \text{const}. \quad (7)$$

$$\mathcal{L}_{\text{U-MAE}}(h) = \mathcal{L}_{\text{MAE}}(h) + \lambda \cdot \mathcal{L}_{\text{unif}}(f), \quad (10)$$

$$\text{where } \mathcal{L}_{\text{unif}}(f) = \mathbb{E}_{x_1} \mathbb{E}_{x_1^-} (f(x_1)^\top f(x_1^-))^2, \quad (11)$$



**Figure 2: Teacher-Student framework** studied in this work. Given an input  $X = [x_1, \dots, x_P]$  (image or text tokens) with  $P$  patches, this framework randomly masks patches to obtain  $\epsilon X = [\epsilon_1 x_1, \dots, \epsilon_p x_P]$  with Bernoulli variable  $\epsilon_p$  to mask, and feeds  $\epsilon X$  into student encoder  $H$  for a latent vector  $H(\epsilon X)$ . Then, student decoder takes  $H(\epsilon X)$  as input and outputs  $h'$  of all patches to predict the output  $h$  of a teacher with vanilla input  $X$  as input. The encoder is two-layer CNN, and the decoder is a linear layer. For **MAE framework**, it has encoder-decoder networks and the decoder have an additional layer to map output of the encoder to recover  $P$  patches (see Fig 6 in Appendix).

# **Test-Time Training**

# TTT-MAE, NeurIPS, 2022

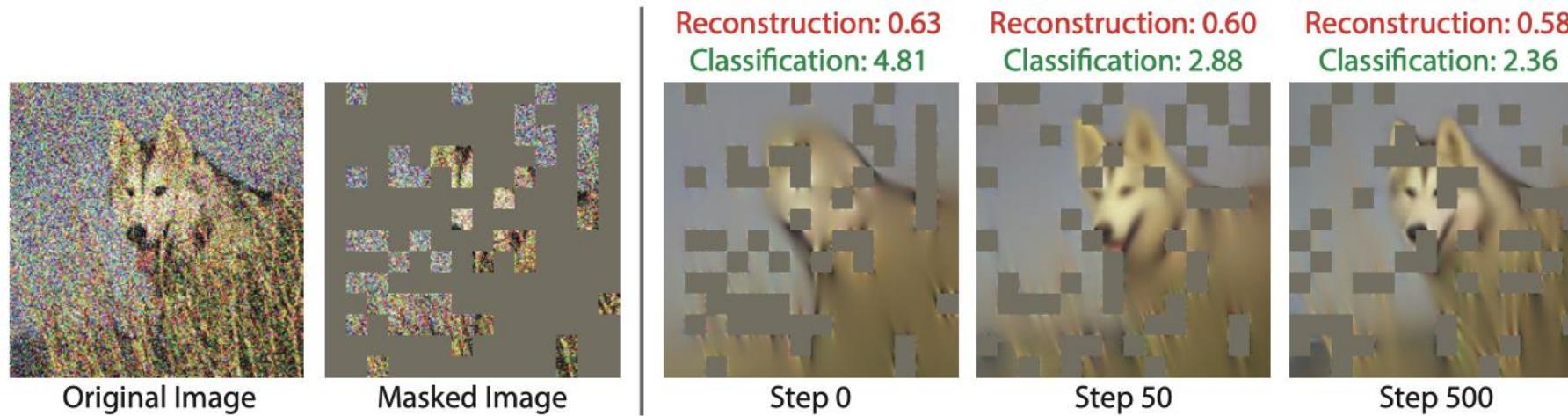


Figure 1: We train an MAE to reconstruct each test image at test time, masking 75% of the input patches. The three reconstructed images on the right visualize the progress of this one-sample learning problem. Loss of the main task (green) – object recognition – keeps dropping even after 500 steps of gradient descent, while the network continues to optimize for reconstruction (red). The unmasked patches are not shown on the right since they are not part of the reconstruction loss.

# MAE as Data Augmentor

# MRA, rejected by NeurIPS 22

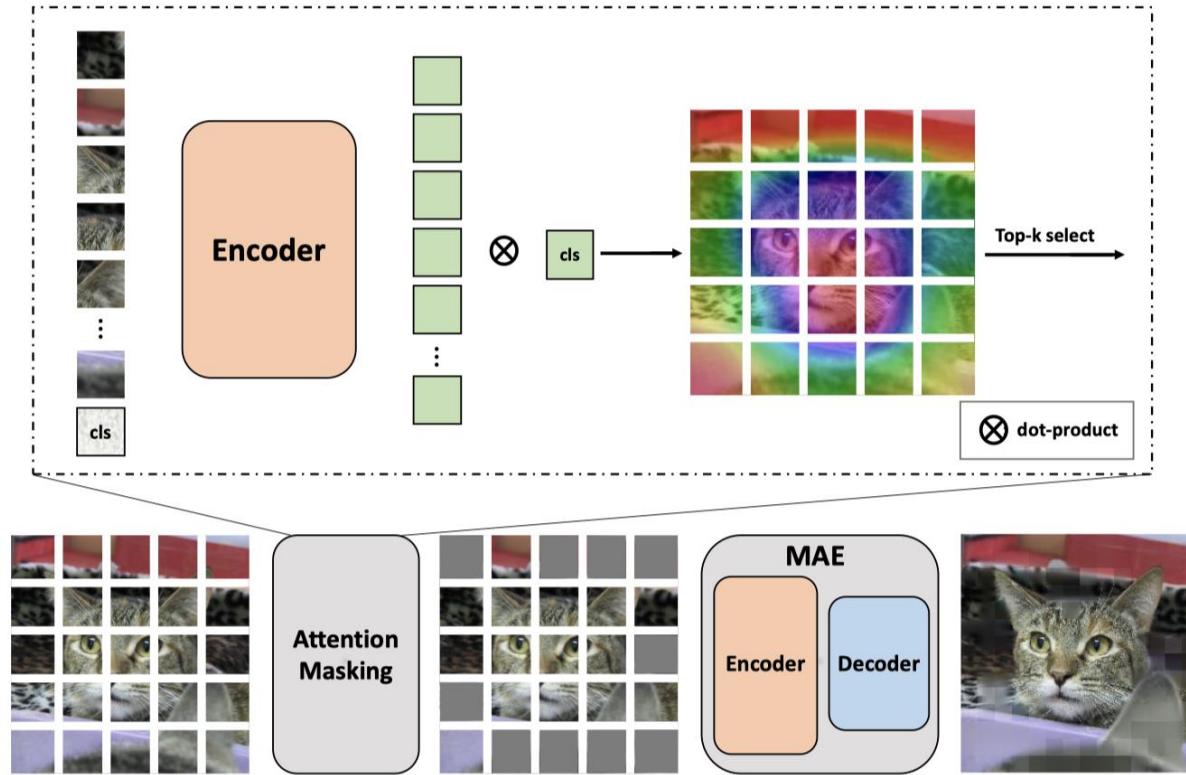
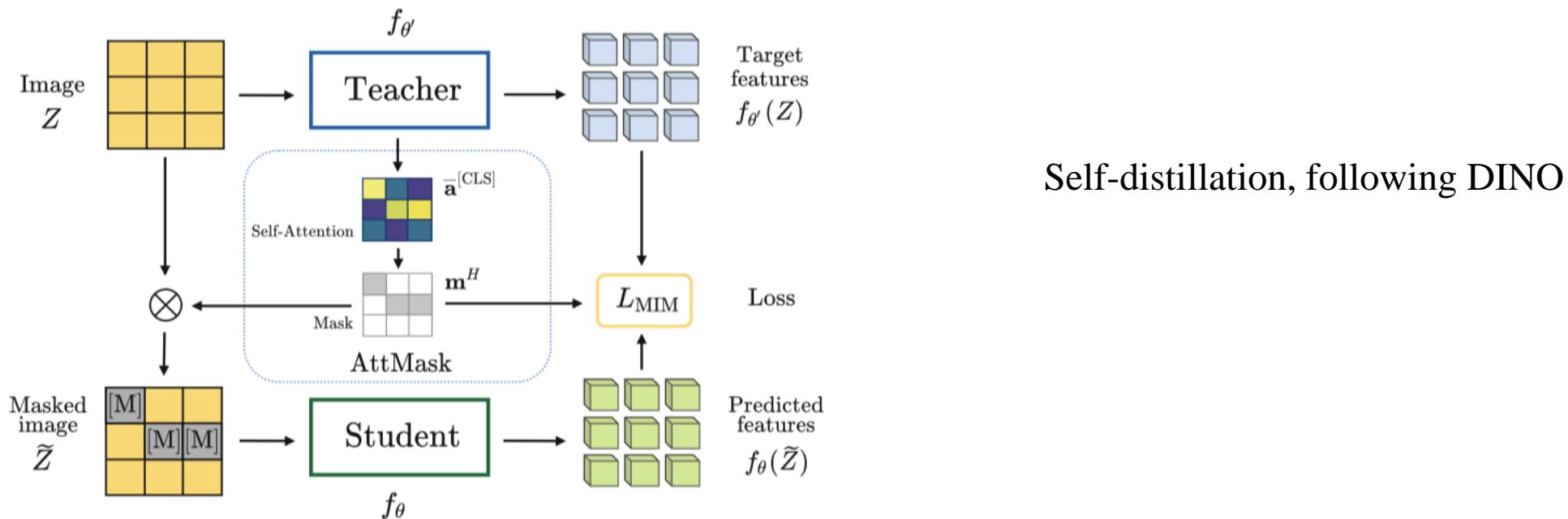


Figure 1: An overview of Mask-Reconstruct Augmentation (MRA). We first mask out the image patches via the attention-based masking strategy. Using class token as query, we calculate attention weights with the keys of each image patch. We remain patches of top-k greatest attention visible. Then, the pretrained masked autoencoder (MAE) [24] completes the original image relying on the visible image patches. The reconstructed image can be viewed as robust augmentation for a number of classification tasks, such as supervised, semi-supervised and few-shot classification.

# SSL Architecture

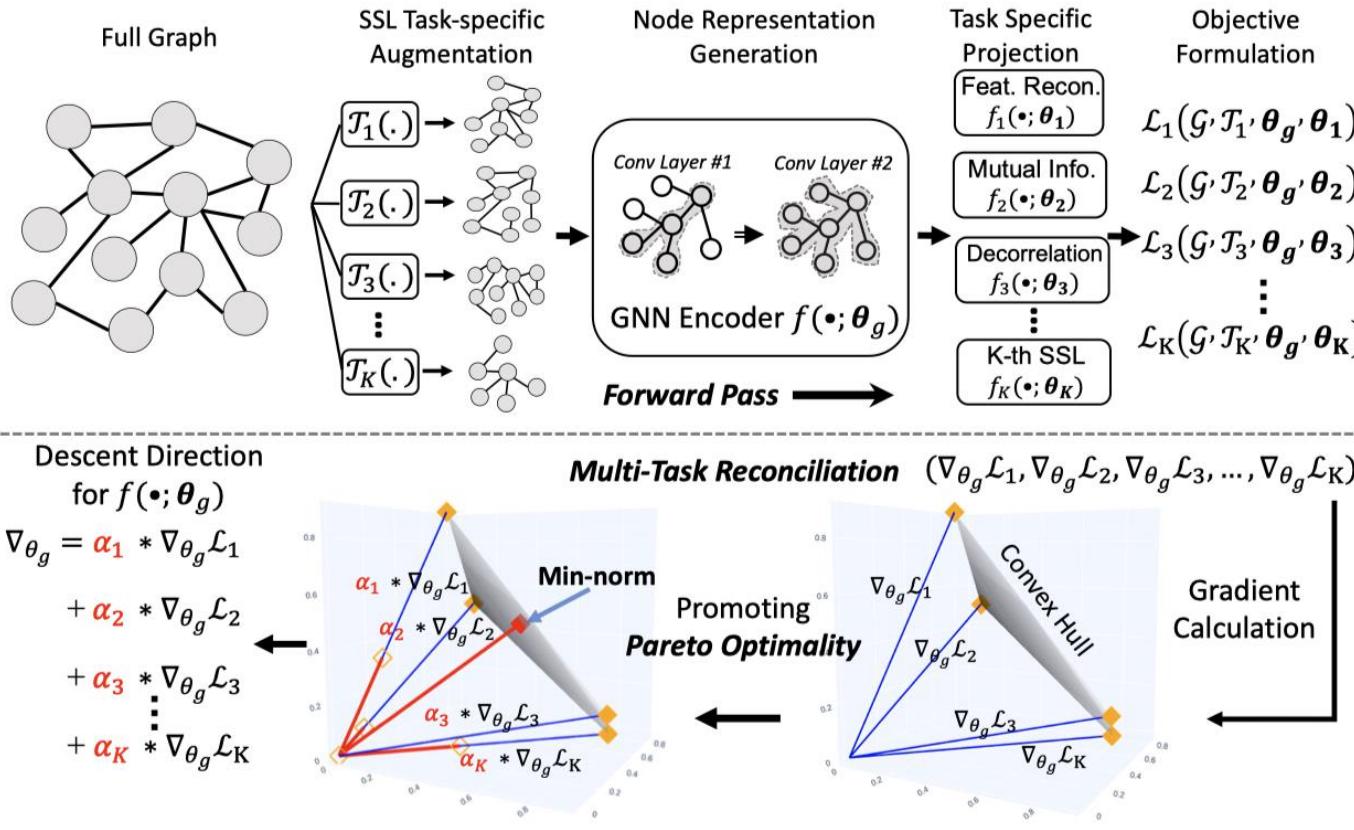
# Self-distillation

AttMask, ECCV, 2022



# Multi-task

ParetoGNN, ICLR, 2023



- **Motivation:**

- Conventional SSL methods only rely on a single philosophy, which could not generalize well on all downstream tasks.
- SSL over multiple pretext tasks leads to better generalization ability and reduces the risk of overfitting.

- **Methods:**

- combine multiple philosophies to enhance task generalization for SSL-based GNNs: **Feature Reconstruction + Topology Reconstruction + whitening decorrelation + mutual information(node-graph) + mutual information(node-subgraph).**
- multiple-gradient descent algorithm to reconcile different philosophies.

# Multi-task

ParetoGNN, ICLR, 2023

Table 2: Performance and task generalization of PARETOGNN as well as the state-of-the-art unsupervised baselines. OOM stands for out-of-memory on a RTX3090 GPU with 24 GB memory.

Method	WIKI.CS	PUBMED	AM.PHOTO	AM.COMP.	Co.CS	Co.PHY.	CHAM.	SQUIRREL	ACTOR	RANK
AVERAGE PERFORMANCE										
DGI	72.50	59.84	83.67	77.91	85.67	79.86	61.34	50.23	33.77	4.4
GRACE	71.01	65.04	80.87	76.06	87.16	OOM	62.12	51.02	32.59	4.8
MVGRL	68.59	63.23	82.49	67.80	82.69	75.72	59.77	45.01	31.22	7.1
AUTOSSL	71.72	65.90	83.96	77.02	85.88	80.03	60.87	49.76	31.33	4.4
BGRL	74.68	67.15	80.14	79.06	87.26	81.42	60.20	49.14	32.33	3.8
CCA-SSG	73.44	64.12	83.62	78.71	83.65	79.77	62.80	51.63	35.90	3.7
GRAPHMAE	67.83	62.15	76.94	72.18	86.80	77.48	58.27	45.22	31.49	6.8
PARETOGNN	<b>76.03</b>	<b>72.48</b>	<b>86.58</b>	<b>82.57</b>	<b>87.80</b>	<b>83.35</b>	<b>65.21</b>	<b>55.31</b>	<b>40.76</b>	<b>1.0</b>
NODE CLASSIFICATION (Accuracy)										
DGI	$75.53 \pm 0.09$	$83.52 \pm 0.52$	$91.61 \pm 0.17$	$83.59 \pm 0.22$	$92.15 \pm 0.25$	$94.51 \pm 0.27$	$61.00 \pm 1.68$	$40.54 \pm 0.62$	$25.19 \pm 0.52$	5.9
GRACE	$80.14 \pm 0.09$	$86.06 \pm 0.35$	$92.78 \pm 0.49$	$89.53 \pm 0.34$	$91.12 \pm 0.25$	OOM	$58.19 \pm 0.53$	$41.36 \pm 0.47$	$24.63 \pm 0.39$	5.0
MVGRL	$79.11 \pm 0.17$	$83.62 \pm 0.10$	$92.48 \pm 0.09$	$82.78 \pm 0.15$	$90.33 \pm 0.17$	$91.05 \pm 0.14$	$49.87 \pm 0.14$	$39.81 \pm 0.40$	$28.01 \pm 0.41$	6.4
AUTOSSL	$79.55 \pm 0.23$	$86.26 \pm 0.20$	$92.71 \pm 0.43$	$88.76 \pm 0.52$	$92.17 \pm 0.17$	$95.13 \pm 0.43$	$58.94 \pm 0.94$	$40.63 \pm 0.62$	$24.54 \pm 0.28$	4.7
BGRL	$82.58 \pm 0.27$	$86.03 \pm 0.17$	$93.17 \pm 0.23$	$90.15 \pm 0.18$	$91.77 \pm 0.47$	$95.73 \pm 0.23$	$56.05 \pm 1.06$	$41.64 \pm 0.79$	$24.03 \pm 0.78$	4.1
CCA-SSG	$82.48 \pm 0.35$	$86.36 \pm 0.35$	$93.29 \pm 0.49$	$89.58 \pm 0.70$	$94.24 \pm 0.17$	$95.63 \pm 0.09$	$57.39 \pm 1.38$	$42.22 \pm 0.94$	$26.35 \pm 0.35$	2.8
GRAPHMAE	$77.12 \pm 0.30$	$83.91 \pm 0.26$	$90.71 \pm 0.40$	$79.44 \pm 0.48$	$93.13 \pm 0.15$	$95.79 \pm 0.06$	$55.50 \pm 0.82$	$35.87 \pm 0.61$	$28.97 \pm 0.27$	5.3
PARETOGNN	<b><math>82.87 \pm 0.13</math></b>	<b><math>87.03 \pm 0.25</math></b>	<b><math>93.85 \pm 0.28</math></b>	<b><math>90.75 \pm 0.17</math></b>	<b><math>92.21 \pm 0.14</math></b>	<b><math>95.45 \pm 0.10</math></b>	<b><math>63.13 \pm 0.84</math></b>	<b><math>46.60 \pm 1.08</math></b>	<b><math>26.62 \pm 0.67</math></b>	<b>1.8</b>
NODE CLUSTERING (NMI)										
DGI	$44.35 \pm 0.12$	$9.68 \pm 0.31$	$60.31 \pm 0.23$	$47.76 \pm 0.02$	$72.88 \pm 0.21$	$58.76 \pm 0.43$	$6.99 \pm 1.74$	$2.16 \pm 0.14$	$1.49 \pm 0.05$	5.1
GRACE	$40.40 \pm 0.10$	$25.64 \pm 0.24$	$55.20 \pm 0.70$	$41.77 \pm 0.32$	<b><math>76.61 \pm 0.26</math></b>	OOM	$11.73 \pm 1.01$	$2.53 \pm 0.31$	$1.09 \pm 0.22$	5.3
MVGRL	$36.20 \pm 0.29$	$24.87 \pm 0.13$	$56.36 \pm 0.08$	$31.27 \pm 0.29$	$73.34 \pm 0.01$	$58.27 \pm 0.01$	<b><math>18.57 \pm 0.26</math></b>	<b><math>4.40 \pm 0.21</math></b>	<b><math>2.63 \pm 0.06</math></b>	4.3
AUTOSSL	$36.99 \pm 0.21$	$28.99 \pm 0.26$	$64.06 \pm 0.65$	$41.85 \pm 0.36$	$74.04 \pm 0.22$	$55.23 \pm 0.18$	$9.67 \pm 1.21$	$2.11 \pm 0.12$	$1.35 \pm 0.11$	5.0
BGRL	$44.95 \pm 0.32$	$26.38 \pm 0.36$	$50.56 \pm 0.24$	$44.04 \pm 0.36$	$74.06 \pm 0.51$	<b><math>61.04 \pm 0.11</math></b>	$11.29 \pm 1.45$	$2.28 \pm 0.53$	$1.32 \pm 0.07$	4.3
CCA-SSG	$44.17 \pm 0.04$	$27.15 \pm 1.56$	$64.06 \pm 0.03$	$49.78 \pm 0.01$	$67.14 \pm 0.49$	$54.33 \pm 2.69$	$12.17 \pm 0.79$	$3.02 \pm 0.04$	$1.01 \pm 0.02$	4.3
GRAPHMAE	$35.73 \pm 0.04$	$19.00 \pm 0.11$	$51.42 \pm 0.46$	$43.51 \pm 0.16$	$76.18 \pm 0.59$	$46.90 \pm 0.07$	$8.66 \pm 0.47$	$3.70 \pm 0.32$	$1.24 \pm 0.05$	5.7
PARETOGNN	<b><math>47.52 \pm 0.29</math></b>	<b><math>34.74 \pm 0.06</math></b>	<b><math>68.25 \pm 1.25</math></b>	<b><math>52.53 \pm 0.34</math></b>	<b><math>74.94 \pm 0.98</math></b>	<b><math>60.43 \pm 0.13</math></b>	<b><math>14.49 \pm 0.23</math></b>	<b><math>2.81 \pm 0.23</math></b>	<b><math>1.53 \pm 0.04</math></b>	<b>1.9</b>
LINK PREDICTION (AUC)										
DGI	$94.36 \pm 0.04$	$78.59 \pm 0.58$	$94.24 \pm 0.16$	$90.37 \pm 0.03$	$93.46 \pm 0.12$	$88.82 \pm 0.05$	$91.64 \pm 0.62$	$92.45 \pm 0.01$	$72.01 \pm 0.41$	5.0
GRACE	$92.32 \pm 0.05$	$87.44 \pm 1.03$	$91.82 \pm 0.12$	$88.58 \pm 0.08$	$97.00 \pm 0.02$	OOM	$93.81 \pm 0.47$	$93.57 \pm 0.06$	$82.31 \pm 0.17$	4.6
MVGRL	$94.34 \pm 0.20$	$90.31 \pm 0.02$	$93.86 \pm 0.25$	$75.15 \pm 0.34$	$92.44 \pm 0.04$	$87.54 \pm 0.08$	$94.50 \pm 0.44$	$88.81 \pm 0.33$	$77.93 \pm 0.07$	5.0
AUTOSSL	$93.86 \pm 0.02$	$86.84 \pm 1.30$	<b><math>95.57 \pm 0.13</math></b>	<b><math>93.99 \pm 0.03</math></b>	$95.71 \pm 0.15$	$95.93 \pm 0.07$	$90.05 \pm 0.78$	<b><math>93.84 \pm 0.19</math></b>	$70.51 \pm 0.29$	4.3
BGRL	$94.31 \pm 0.06$	$94.35 \pm 0.02$	$93.33 \pm 0.06$	$93.59 \pm 0.05$	$97.37 \pm 0.04$	$93.38 \pm 0.07$	$92.24 \pm 0.62$	$83.60 \pm 0.10$	$76.91 \pm 0.10$	4.4
CCA-SSG	$90.92 \pm 0.33$	$73.53 \pm 1.50$	$89.47 \pm 0.13$	$86.72 \pm 0.13$	$91.77 \pm 0.03$	$93.64 \pm 0.02$	<b><math>95.43 \pm 0.31</math></b>	$87.30 \pm 0.25$	$78.14 \pm 0.09$	5.7
GRAPHMAE	$89.47 \pm 0.02$	$86.24 \pm 0.19$	$83.33 \pm 0.07$	$84.65 \pm 1.08$	$96.48 \pm 0.17$	<b><math>96.57 \pm 0.08</math></b>	$89.40 \pm 1.48$	$86.48 \pm 0.14$	$80.65 \pm 0.28$	5.9
PARETOGNN	<b><math>96.48 \pm 0.01</math></b>	<b><math>94.58 \pm 0.02</math></b>	<b><math>96.08 \pm 0.08</math></b>	<b><math>97.16 \pm 0.04</math></b>	<b><math>98.18 \pm 0.01</math></b>	<b><math>98.33 \pm 0.12</math></b>	<b><math>95.78 \pm 0.05</math></b>	<b><math>96.46 \pm 0.05</math></b>	<b><math>84.29 \pm 0.04</math></b>	<b>1.0</b>
PARTITION PREDICTION (Accuracy)										
DGI	$75.75 \pm 0.19$	$67.58 \pm 0.01$	<b><math>88.50 \pm 0.20</math></b>	<b><math>89.91 \pm 0.07</math></b>	$84.19 \pm 0.66$	<b><math>77.34 \pm 0.01</math></b>	$85.75 \pm 0.14$	$65.77 \pm 0.02$	$36.40 \pm 0.02$	2.9
GRACE	$71.19 \pm 0.13$	$61.01 \pm 1.57$	$83.66 \pm 0.83$	$84.36 \pm 0.21$	$83.89 \pm 0.82$	OOM	$84.76 \pm 2.03$	$66.64 \pm 0.12$	$22.32 \pm 0.05$	5.4
MVGRL	$64.73 \pm 0.25$	$54.11 \pm 0.06$	$87.25 \pm 0.44$	$81.99 \pm 0.04$	$74.66 \pm 0.19$	$66.03 \pm 0.35$	$76.12 \pm 1.03$	$47.03 \pm 0.81$	$16.31 \pm 0.04$	7.2
AUTOSSL	$76.47 \pm 0.07$	$61.52 \pm 0.01$	$83.51 \pm 0.05$	$83.46 \pm 0.79$	$81.59 \pm 0.80$	$73.85 \pm 0.04$	$84.83 \pm 0.94$	$62.45 \pm 0.02$	$28.90 \pm 0.04$	5.0
BGRL	$76.87 \pm 0.03$	$61.84 \pm 0.02$	$83.51 \pm 0.05$	$88.45 \pm 0.38$	$85.84 \pm 1.29$	$75.53 \pm 0.01$	$81.21 \pm 0.99$	$69.04 \pm 0.39$	$27.07 \pm 0.04$	3.9
CCA-SSG	$76.18 \pm 0.46$	$69.46 \pm 0.26$	$87.67 \pm 0.31$	$88.77 \pm 0.19$	$81.44 \pm 0.23$	$75.49 \pm 0.15$	$86.22 \pm 0.66$	$73.97 \pm 0.22$	$38.11 \pm 0.29$	3.1
GRAPHMAE	$69.01 \pm 0.27$	$59.46 \pm 0.08$	$82.31 \pm 0.13$	$81.11 \pm 1.57$	$81.39 \pm 0.20$	$70.67 \pm 0.16$	$79.50 \pm 1.26$	$54.82 \pm 0.67$	$15.09 \pm 0.01$	7.2
PARETOGNN	<b><math>77.23 \pm 0.36</math></b>	<b><math>73.57 \pm 0.23</math></b>	$88.13 \pm 0.39$	$89.84 \pm 0.06$	<b><math>85.89 \pm 0.33</math></b>	<b><math>79.19 \pm 0.20</math></b>	<b><math>87.43 \pm 1.05</math></b>	<b><math>75.39 \pm 0.65</math></b>	<b><math>50.61 \pm 0.75</math></b>	<b>1.2</b>

# Multi-task

AutoSSL, ICLR, 2022

## AUTOMATED SELF-SUPERVISED LEARNING FOR GRAPHS

**Wei Jin** \*

Michigan State University  
jinwei2@msu.edu

**Xiaorui Liu**

Michigan State University  
xiaorui@msu.edu

**Xiaoyu Zhao**

City University of Hong Kong  
xy.zhao@cityu.edu.hk

**Yao Ma**

New Jersey Institute of Technology  
yao.ma@njit.edu

**Neil Shah**

Snap Inc.  
nshah@snap.com

**Jiliang Tang**

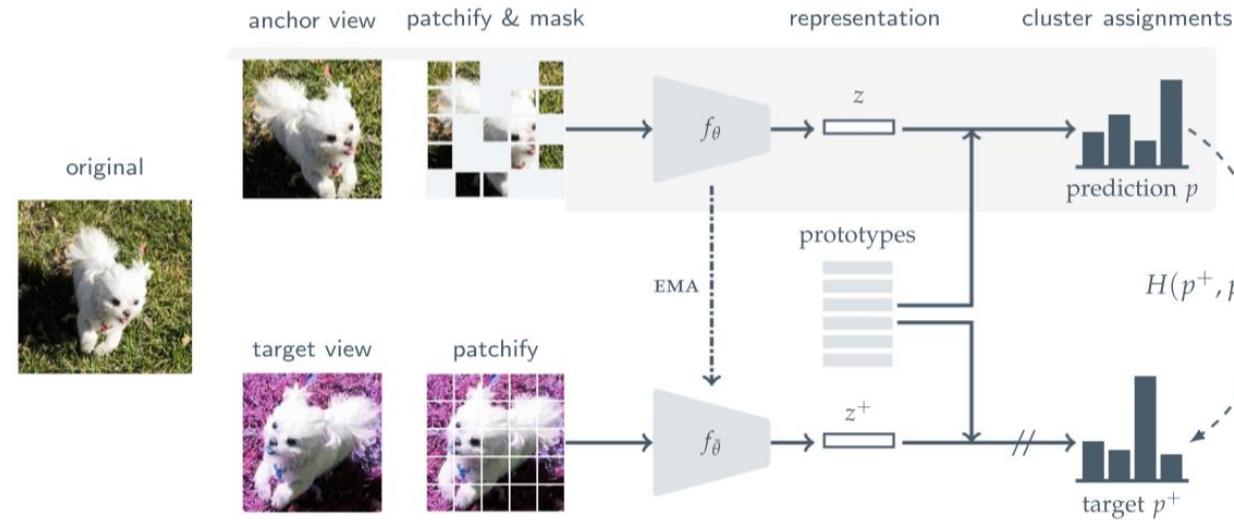
Michigan State University  
tangjili@msu.edu

### ABSTRACT

Graph self-supervised learning has gained increasing attention due to its capacity to learn expressive node representations. Many pretext tasks, or loss functions have been designed from distinct perspectives. However, we observe that different pretext tasks affect downstream tasks differently across datasets, which suggests that searching over pretext tasks is crucial for graph self-supervised learning. Different from existing works focusing on designing single pretext tasks, this work aims to investigate how to automatically leverage multiple pretext tasks effectively. Nevertheless, evaluating representations derived from multiple pretext tasks without direct access to ground truth labels makes this problem challenging. To address this obstacle, we make use of a key principle of many real-world graphs, i.e., homophily, or the principle that “like attracts like,” as the guidance to effectively search various self-supervised pretext tasks. We provide theoretical understanding and empirical evidence to justify the flexibility of homophily in this search task. Then we propose the AUTOSSL framework to automatically search over combinations of various self-supervised tasks. By evaluating the framework on 8 real-world datasets, our experimental results show that AUTOSSL can significantly boost the performance on downstream tasks including node clustering and node classification compared with training under individual tasks.

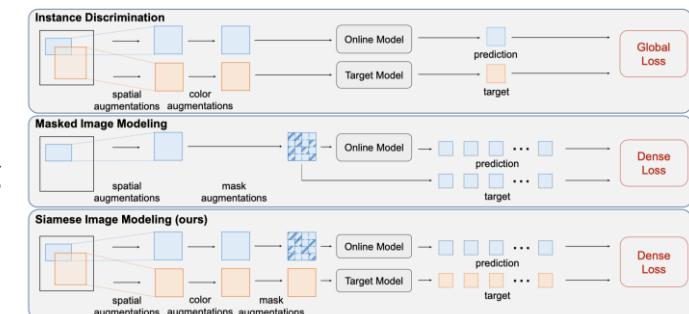
# Masked Siamese Networks for Label-Efficient Fine-Tuning

MSN, ECCV, 2022



Related work (Siamese Image Modeling):

[1] Chenxin Tao, Xizhou Zhu, Weijie Su, Gao Huang, Bin Li, Jie Zhou, Yu Qiao, Xiaogang Wang, Jifeng Dai. "Siamese Image Modeling for Self-Supervised Vision Representation Learning." arXiv. 2022



[1] Ioannis Kakogeorgiou, Spyros Gidaris, Bill Psomas, Yannis Avrithis, Andrei Bursuc, Konstantinos Karantzalos, Nikos Komodakis. "What to Hide from Your Students: Attention-Guided Masked Image Modeling." In ECCV, 2022

# Maximizing Rate Reduction

G<sup>2</sup>R, WWW, 2023

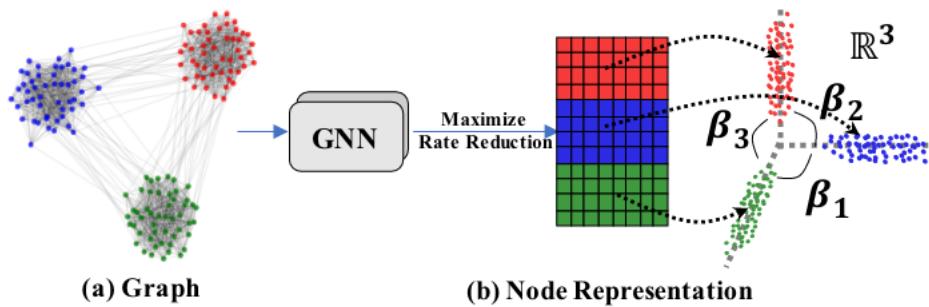


Figure 1: Overview of G<sup>2</sup>R. It maps nodes in distinct groups (implicitly preserved in adjacency matrix) into different subspaces, while each subspace is compact and different subspaces are dispersedly distributed. Different colors indicate different subspaces and node groups.

# Masking Strategy

Masking mechanisms define the specific pretext task, i.e., what kind of information is to be exploited and what kind of information is to be predicted.

What to Mask?

- Random Masking
- Adversarial-based Masking
- Attention-based Masking
- Semantic-based Masking

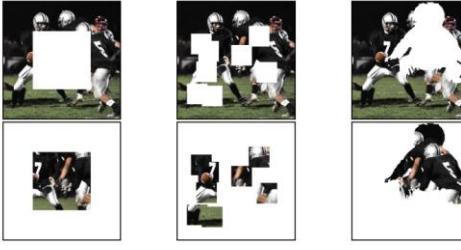
How to Mask?

- Zeros-replacement, Noise-replacement, other-feature-replacement, learned shared token
- Corruption-based
- Mix-based

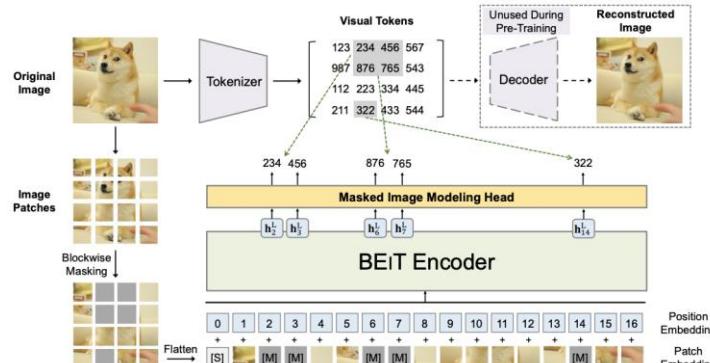
# What to Mask?

- Random Masking
- Adversarial-based Masking
- Attention-based Masking
- Semantic-based Masking

# Masking Strategies in MIM - Random



Context Encoder, CVPR16  
Central region mask, random shape mask



BEiT, ICLR22  
Random block mask

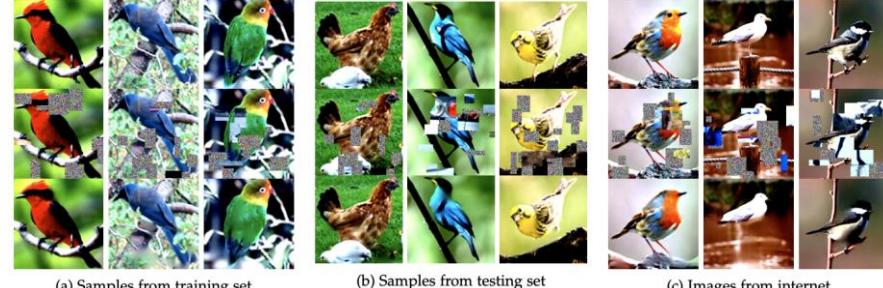
[1] Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, Alexei A. Efros. “Context Encoders: Feature Learning by Inpainting.” In *CVPR*, 2016

[2] Sara Atito. “SiT: Self-Supervised VIision Transformer.” arXiv. 2021

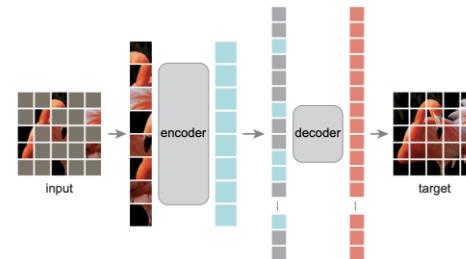
[3] Li Dong, Songhao Piao, Furu Wei. “BEiT: BERT Pre-Training of Image Transformers.” In *ICLR*, 2022

[4] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, Ross Girshick. “Masked Autoencoders Are Scalable Vision Learners.” In *CVPR*, 2022

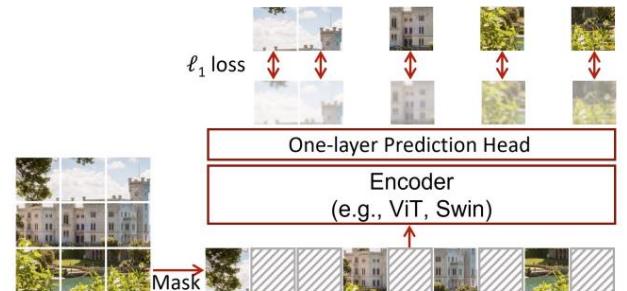
[5] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, Han Hu. “SimMIM: A Simple Framework for Masked Image Modeling.” In *CVPR*, 2022



SiT, arXiv21  
Random block mask



MAE, CVPR22  
Random patch mask



SimMIM, CVPR22  
Random patch mask

# Masking Strategies in MIM - Adversarial

## Masked Language Model

BERT    A bird with a small head, yellow belly and short tail. → A bird with a small [pink box], yellow [brown box] and short [green box].

## Masked Image Models

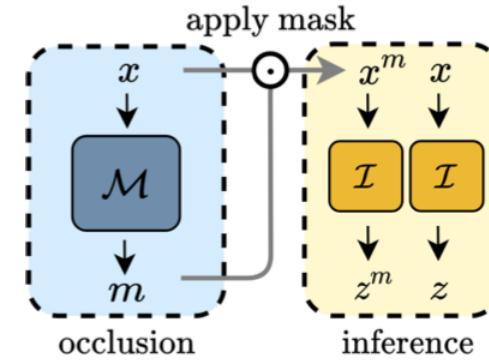


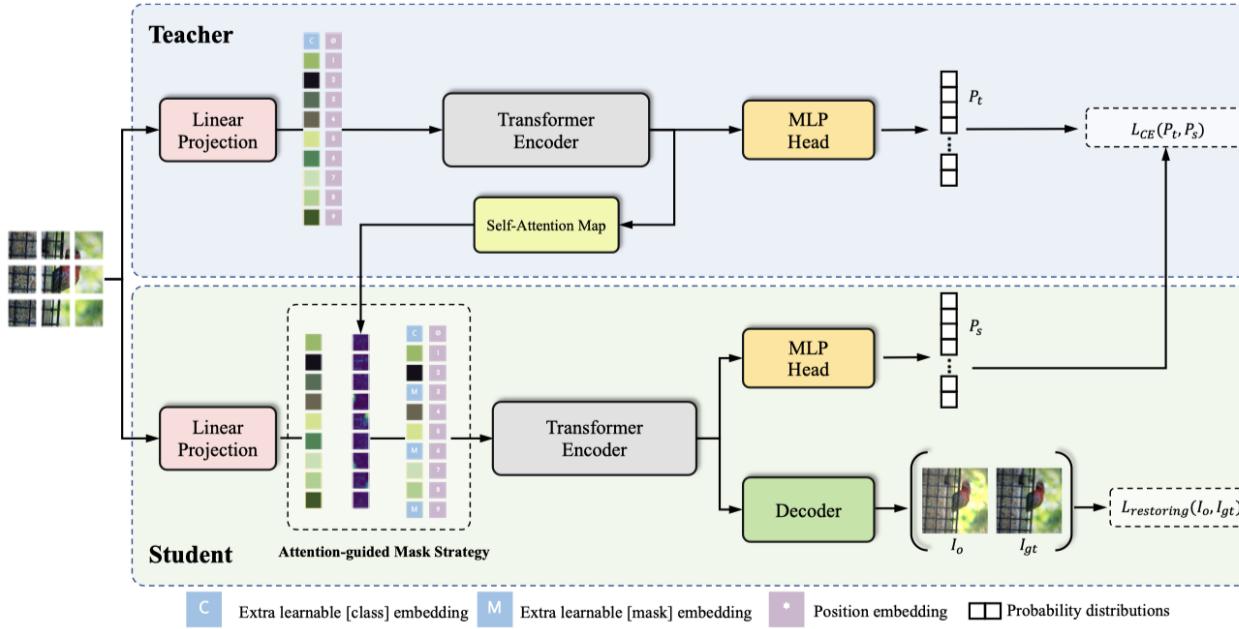
Figure 2. ADIOS Architecture.

$$\mathcal{I}^*, \mathcal{M}^* = \arg \min_{\mathcal{I}} \max_{\mathcal{M}} \mathcal{L}(\mathbf{x}; \mathcal{I}, \mathcal{M}). \quad (1)$$

$$\mathcal{L}_{AE}(\mathbf{x}; \mathcal{I}, \mathcal{M}) = \mathcal{D}(\mathbf{x}, \hat{\mathbf{x}}) = \mathcal{D}(\mathbf{x}, \mathcal{I}(\mathbf{x} \odot \mathcal{M}(\mathbf{x}))), \quad (2)$$

ADIOS, ICML22  
“learn to mask” by adversarial training

# Masking Strategies in MIM - Attention



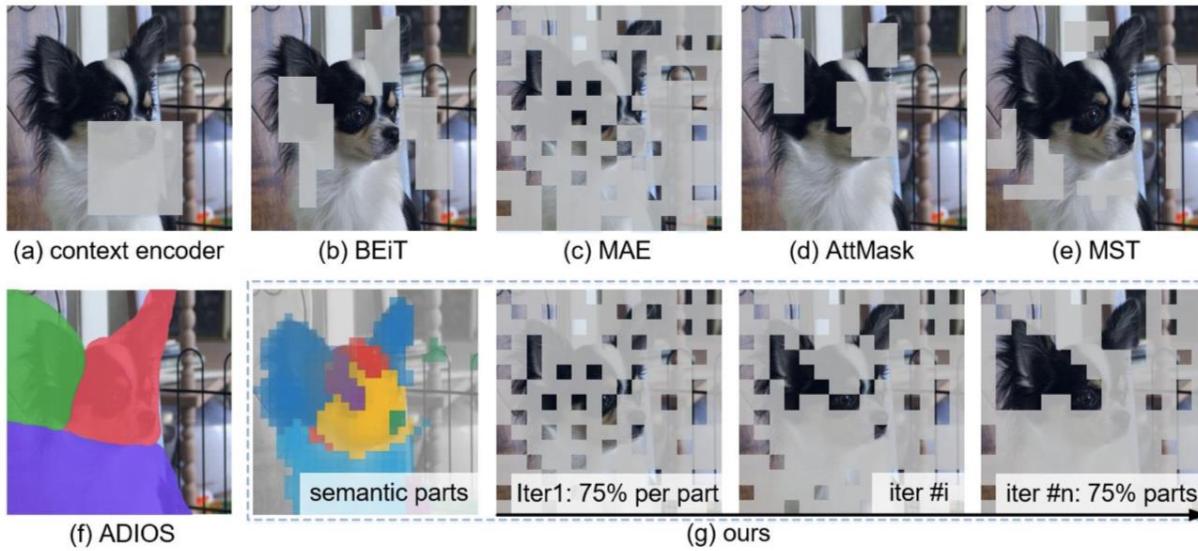
MST, NeurIPS21  
Attention-guided mask



AttMask, ECCV22  
Attention-guided mask

- [1] Zhaowen Li, Zhiyang Chen, Fan Yang, Wei Li, Yousong Zhu, Chaoyang Zhao, Rui Deng, et al. “MST: Masked Self-Supervised Transformer for Visual Representation.” In *NeurIPS*. 2021
- [2] Ioannis Kakogeorgiou, Spyros Gidaris, Bill Psomas, Yannis Avrithis, Andrei Bursuc, Konstantinos Karantzalos, Nikos Komodakis. “What to Hide from Your Students: Attention-Guided Masked Image Modeling.” In *ECCV*, 2022

# Masking Strategies in MIM - Semantic



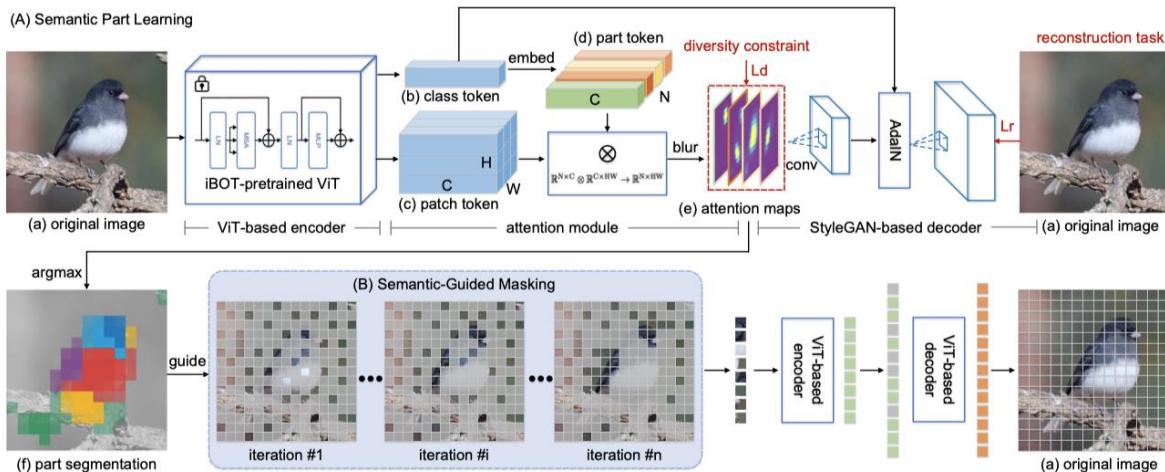
SemMAE, NeurIPS22

## Motivation:

There is still a large gap for masked autoencoding (MAE) between vision and language due to different signal natures. A sentence can be semantically decomposed into words, while the semantic decomposition of an image is not trivial to be obtained.

## Methods:

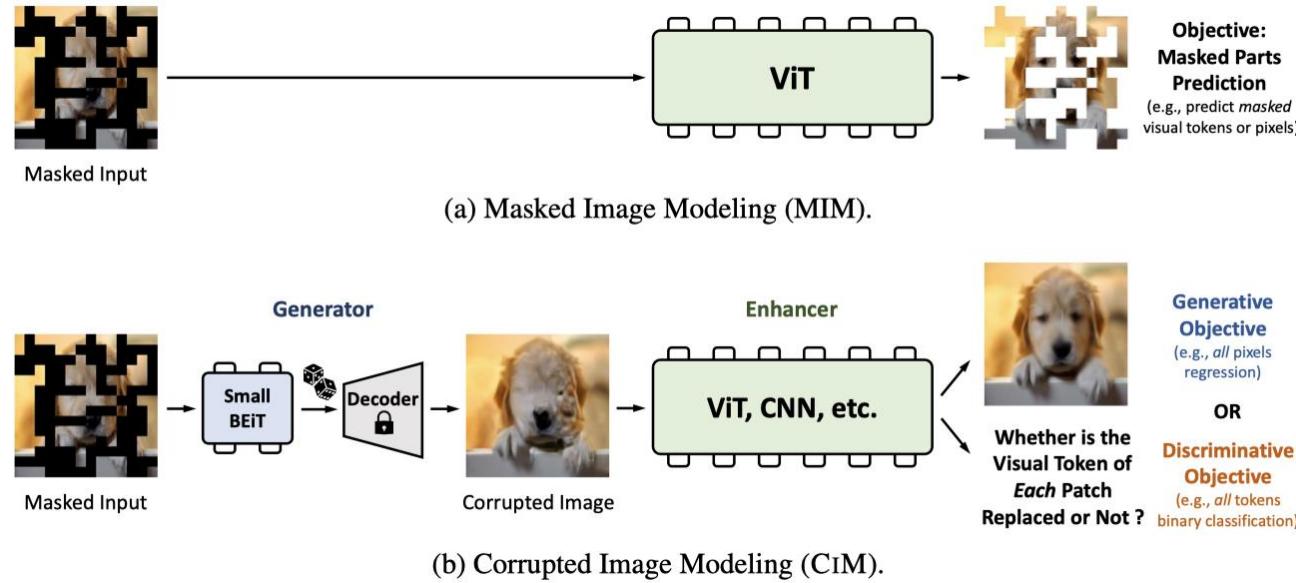
1. Identify semantic parts.
2. Dynamic easy-to-hard masking strategies:
  - a. Mask a portion of patches in each part, to learn intra-part patterns
  - b. Mask a portion of (whole) parts in an image, to learn inter-part patterns



[1] Gang Li, Heliang Zheng, Daqing Liu, Chaoyue Wang, Bing Su, Changwen Zheng. "SemMAE: Semantic-Guided Masking for Learning Masked Autoencoders." In *NeurIPS*. 2022

# How to Mask?

# CIM, ICLR, 2023



- **Motivation:**

- **Deficiencies of [MASK] token:**
  - MIM is tightly coupled with the Transformer family, and cannot be combined with CNN family.
  - It causes a discrepancy between pre-training and fine-tuning.

- **Methods:**

- Use corrupt-reconstruct to replace mask-reconstruct
- Use a small trainable BEiT as an generator to corrupt the input image.

# CIM, ICLR, 2023

Table 1: ImageNet-1K end-to-end fine-tuning top-1 accuracy of vanilla ViT-Small/16 and ViT-Base/16 models.

<sup>†</sup>Doubled attention heads. <sup>‡</sup>Our reproduction.

Models	PT Epochs	Top-1
<i>ViT-Small/16 model results</i>		
Scratch ( <a href="#">Touvron et al., 2021a</a> )		79.9
MoCo-v3 <sup>†</sup> ( <a href="#">Chen et al., 2021</a> )	600	81.4
DINO ( <a href="#">Caron et al., 2021</a> )	1600	81.5
BEiT ( <a href="#">Bao et al., 2021</a> )	300	81.3
<b>CIM-RESPIX (Ours)</b>	300	81.5
<b>CIM-REVDET (Ours)</b>	300	<b>81.6</b>
<i>ViT-Base/16 model results</i>		
Scratch ( <a href="#">Touvron et al., 2021a</a> )		81.8
Scratch ( <a href="#">He et al., 2021</a> )		82.3
DINO ( <a href="#">Caron et al., 2021</a> )	1600	82.8
MoCo-v3 ( <a href="#">Chen et al., 2021</a> )	600	83.2
BEiT ( <a href="#">Bao et al., 2021</a> )	300	82.9
BEiT ( <a href="#">Bao et al., 2021</a> )	800	83.2
MAE <sup>‡</sup> ( <a href="#">He et al., 2021</a> )	800	83.1
<b>CIM-REVDET (Ours)</b>	300	<b>83.3</b>
<b>CIM-RESPIX (Ours)</b>	300	<b>83.3</b>

# MixMask, arXiv, 2022

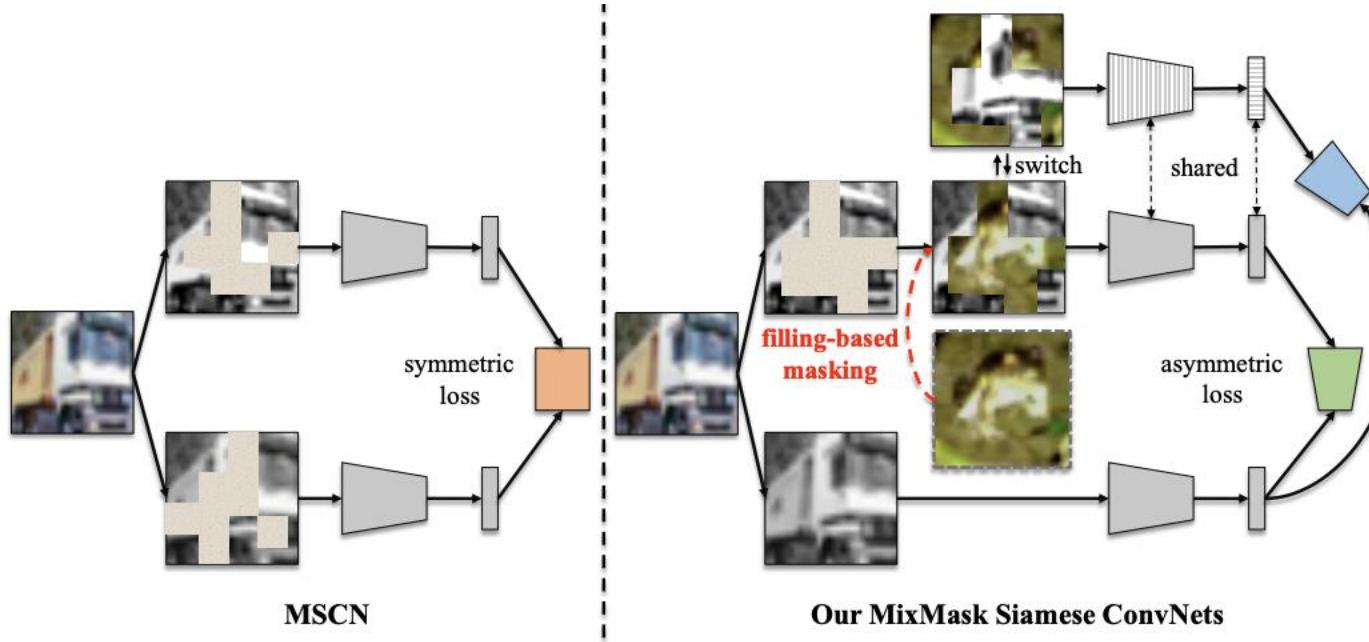


Figure 2. Illustration of the Masked Siamese ConvNets (left) and our proposed framework (right). MixMask branch incorporates asymmetry into the loss function design by generating images with different rates of similarity to the images in the original branch. In MixMask branch image of the truck is presented twice with different levels of similarity to the image in the original branch due to the regions masked with contents of another image.

## Future Directions

- Adversarial-based
- Attention-guided (combined with GT)
- Semantic-guided
  - What is semantic part in graph? Subgraph partition? Saliency? Prototype?
- Heuristic:
  - Structure-guided
    - Degree/Homophily-guided
    - Ego-graph-based
  - Skip-mask
  - Dynamic mask: easy-to-hard, fine-to-coarse
  - Multi-granularity
- Combine masking and augmentation